

# DETECTING PHYLOGENETIC RELATIONS OUT FROM SPARSE CONTEXT TREES

FLORENCIA LEONARDI, SERGIO R. MATIOLI, HUGO A. ARMELIN,  
AND ANTONIO GALVES

**ABSTRACT.** The goal of this paper is to study the similarity between sequences using a distance between the *context* trees associated to the sequences. These trees are defined in the framework of *Sparse Probabilistic Suffix Trees* (SPST), and can be estimated using the SPST algorithm. We implement the Phyl-SPST package to compute the distance between the sparse context trees estimated with the SPST algorithm. The distance takes into account the structure of the trees, and indirectly the transition probabilities. We apply this approach to reconstruct a phylogenetic tree of protein sequences in the globin family of vertebrates. We compare this tree with the one obtained using the well-known PAM distance.

## 1. INTRODUCTION

In this work we propose to use the framework of Sparse Probabilistic Suffix Trees (SPST) to analyze the similarity between sequences and to infer the evolution of protein families. SPST was first introduced in Leonardi and Galves (2005) as a generalization of the PST algorithm, proposed in Ron et al. (1996). SPST has shown to be useful in protein modeling and classification, performing better than the PST algorithm (Leonardi; 2006). The model that inspired the SPST algorithm is a generalization of Variable Length Markov Chains (VLMC), introduced by Rissanen (1983), and takes into account the property of sparseness of the sequences. Given a sequence, SPST estimates a set of *sparse contexts*. A sparse context is a short sequence of sub-sets of symbols (in a given alfabet) that are relevant to predict any symbol in the sequence, given that the preceding symbols belong to the sub-sets of the context. The SPST algorithm also estimates the transition probabilities associated to each context. The transition probabilities give the probability of each symbol conditioned on the fact that the preceding symbols belong to the sparse context.

An interesting property of the set of sparse contexts is that it induces a partition of the set of all possible sequences and can be represented as a tree. We use this partition property to define a distance between context trees. This distance can be used to measure the similarity between protein sequences.

To our knowledge it has not been proposed yet in the literature a method for sequence comparison using the information contained in the architecture of the context trees associated to the sequences. The more closely related approaches proposed until date are those that model the sequences as first order Markov chains and use a statistical measure to infer the similarity between them (Wu et al.; 2001; Pham and Zuegg; 2004). The more remarkable difference between these approaches and our is that we do not use directly the estimated probabilities of the model. Instead of that we use the context tree architecture, that is trivial in first order Markov chains. We show here that the context tree architecture can have important structural information that may be useful to measure the similarity between sequences.

The paper is organized as follows. In Section 2 we review some definitions in the framework of SPST. In Section 3 we introduce the distance between sparse trees. In Section 4 we present the results obtained for the globin protein family of vertebrates and finally in Section 5 we discuss some aspects of our method.

## 2. SPARSE CONTEXT TREES

Let  $A$  be a finite alphabet (for example, the set of twenty amino acids) of size  $|A|$ . We will denote by  $\mathcal{P}_A$  the set of parts of  $A$ . That is,

$$\mathcal{P}_A = \{v : v \subset A\}.$$

The elements in  $\mathcal{P}_A^j$  will be denoted by  $w = (w_{-j}, \dots, w_{-1})$ . On the other hand, we will denote by  $\mathcal{P}_A^*$  the set of all finite sequences of elements in  $\mathcal{P}_A$ ; that is,

$$\mathcal{P}_A^* = \bigcup_{j=1}^{\infty} \mathcal{P}_A^j.$$

**Definition 2.1.** Let  $(X_t)_{t \in \mathbb{N}}$  be a stochastic process taking values on the finite alphabet  $A$ . We will say that the process  $(X_t)_{t \in \mathbb{N}}$  is a *sparse stochastic chain* if there exists a set  $\tau \subset \mathcal{P}_A^*$  such that:

- (1) For any sequence  $x_0, \dots, x_n$  satisfying

$$\mathbb{P}[X_0 = x_0, \dots, X_{n-1} = x_{n-1}] > 0,$$

there exists an element  $(w_{-k}, \dots, w_{-1}) \in \tau$  such that

$$\begin{aligned} \mathbb{P}[X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_0 = x_0] = \\ \mathbb{P}[X_n = x_n | X_{n-1} \in w_{-1}, \dots, X_{n-k} \in w_{-k}]. \end{aligned} \quad (2.2)$$

- (2) If  $(w_{-k}, \dots, w_{-1})$  and  $(\bar{w}_{-\bar{k}}, \dots, \bar{w}_{-1})$  belong to  $\tau$  and there exists  $j$  such that  $w_{-i} \cap \bar{w}_{-i} \neq \emptyset$  for  $i = 1, \dots, j$ , then  $w_{-i} = \bar{w}_{-i}$  for  $i = 1, \dots, j$ .

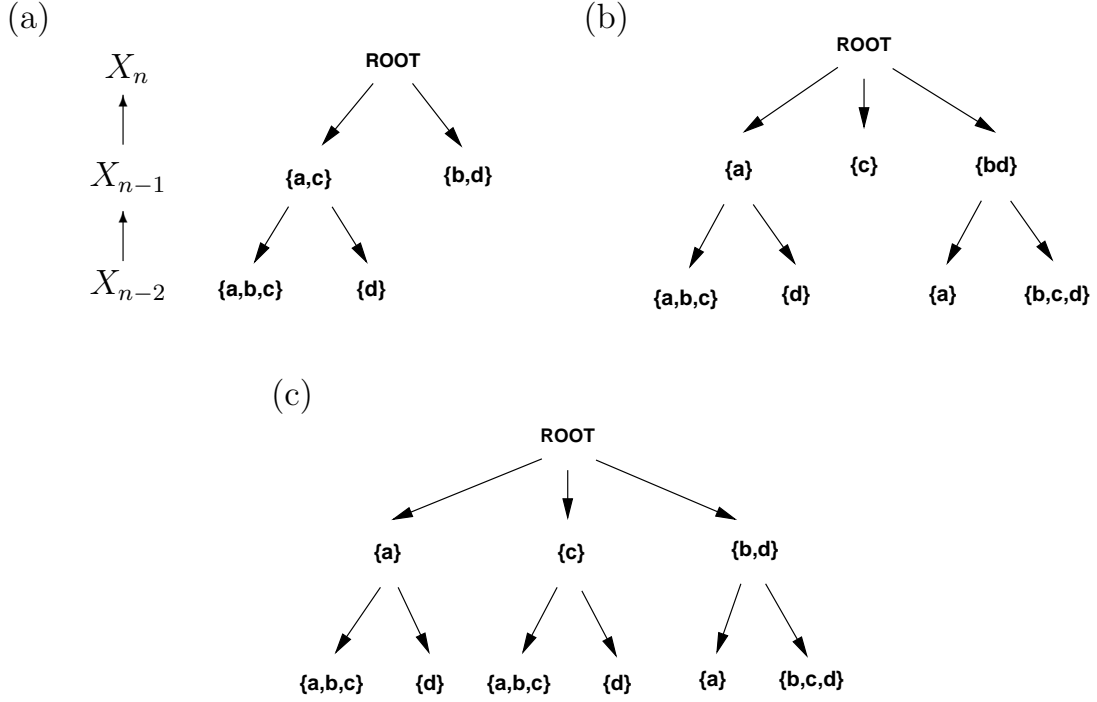


FIGURE 1. Examples of sparse trees over the alphabet  $A = \{a, b, c, d\}$ .

(a) The index of the variables grows in the direction from the leaves to the root. In this case, the set of sparse contexts is  $\{(\{a, b, c\}, \{a, c\}), (\{d\}, \{a, c\}), (\{b, d\})\}$ . (c) Maximum between the trees in (a) and (b).

- (3) The set  $\tau$  is the *minimum* that satisfies 1. and 2. That is; if  $\bar{\tau}$  satisfies 1. and 2. then, for any  $(\bar{w}_{-\bar{k}}, \dots, \bar{w}_{-1}) \in \bar{\tau}$  there exists  $(w_{-k}, \dots, w_{-1}) \in \tau$  such that  $\bar{k} \geq k$  and  $\bar{w}_j \subset w_j$  for all  $j = 1, \dots, k$ .

Each sequence  $(w_{-k}, \dots, w_{-1}) \in \tau$  is called *sparse context* and the set  $\tau$  is called *sparse context tree*. This name is justified because the set of sparse contexts can be represented as a rooted tree. In this tree, each context  $w = (w_{-k}, \dots, w_{-1})$  is represented by a complete branch, in which the first node on top is  $w_{-1}$  and so on until the last element  $w_{-k}$  which is represented by the terminal node of the branch (Fig. 1).

Recently, it was proposed an algorithm to estimate the set of sparse contexts and the transition probabilities given by 2.2 (Leonardi and Galves; 2005; Leonardi; 2006). This algorithm represents internally the set of sparse contexts as a tree, as described above. We believe that this tree contains important structural information that can be used to measure the similarity between sequences. Our goal in this paper is to

show some results concerning this conjecture. With this aim we propose to use a distance between sparse context trees to measure the relatedness between symbolic sequences. This distance is defined in the next section.

### 3. A METRIC SPACE OF SPARSE TREES

Given a sparse context  $w = (w_{-k}, \dots, w_{-1})$  we denote by  $l(w)$  its length, that is  $l(w) = k$ . We use the notation  $s(w)$  for the product of the cardinals of the  $w_i$ 's, that is

$$s(w) = \prod_{i=1}^{l(w)} |w_i|,$$

where  $|w_i|$  is the number of symbols in  $w_i$ .

Given two sparse contexts  $w = (w_{-k}, \dots, w_{-1})$  and  $\bar{w} = (\bar{w}_{-\bar{k}}, \dots, \bar{w}_{-1})$  we define the intersection between  $w$  and  $\bar{w}$  (assuming without loss of generality that  $k \geq \bar{k}$ ) by  $w \cap \bar{w} = (w_{-k}, \dots, w_{-(\bar{k}+1)}, w_{-\bar{k}} \cap \bar{w}_{-\bar{k}}, \dots, w_{-1} \cap \bar{w}_{-1})$ , if  $w_i \cap \bar{w}_i \neq \emptyset$  for all  $i = 1, \dots, \bar{k}$ . In the case  $w_i \cap \bar{w}_i = \emptyset$  for some  $i = 1, \dots, \bar{k}$  we define  $w \cap \bar{w} = \emptyset$ .

Given two sparse trees  $\tau = \{w^1, \dots, w^n\}$  and  $\bar{\tau} = \{\bar{w}^1, \dots, \bar{w}^m\}$ , we define the maximum between  $\tau$  and  $\bar{\tau}$  by

$$\tau \vee \bar{\tau} = \{w^i \cap \bar{w}^j \mid w^i \cap \bar{w}^j \neq \emptyset; i = 1, \dots, n; j = 1, \dots, m\}.$$

The maximum between the trees of Figure 1(a)-(b) can be seen in Figure 1(c).

Before defining the distance between sparse context trees we introduce the notion of  $\beta$ -entropy of a tree  $\tau$ . Following Simovici and Szymon (2006) we define, for all  $\beta > 0$ ,

$$\mathcal{H}_\beta(\tau) = \frac{1}{2^{1-\beta} - 1} \left( \sum_{w \in \tau} [s(w) |A|^{-l(w)}]^\beta - 1 \right), \quad \text{if } \beta \neq 1,$$

and

$$\mathcal{H}_\beta(\tau) = - \sum_{w \in \tau} s(w) |A|^{-l(w)} \cdot \log_2 [s(w) |A|^{-l(w)}], \quad \text{if } \beta = 1.$$

Then, given two sparse trees,  $\tau$  and  $\bar{\tau}$ , we define the  $\beta$ -distance between them as

$$d_\beta(\tau, \bar{\tau}) = 2\mathcal{H}_\beta(\tau \wedge \bar{\tau}) - \mathcal{H}_\beta(\tau) - \mathcal{H}_\beta(\bar{\tau}). \quad (3.1)$$

It can be seen that  $d_\beta(\cdot, \cdot)$  defines a distance over the set of all context trees. The proof of this assertion can be found in Simovici and Szymon (2006).

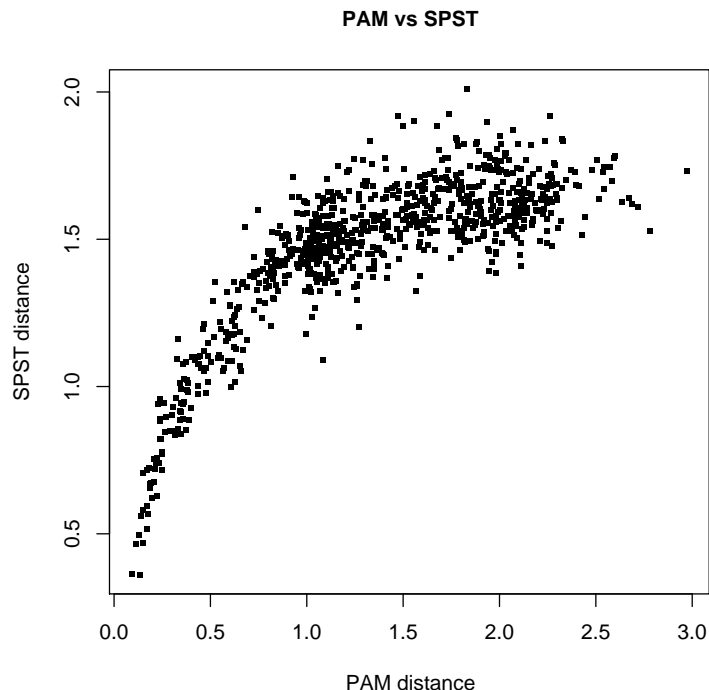


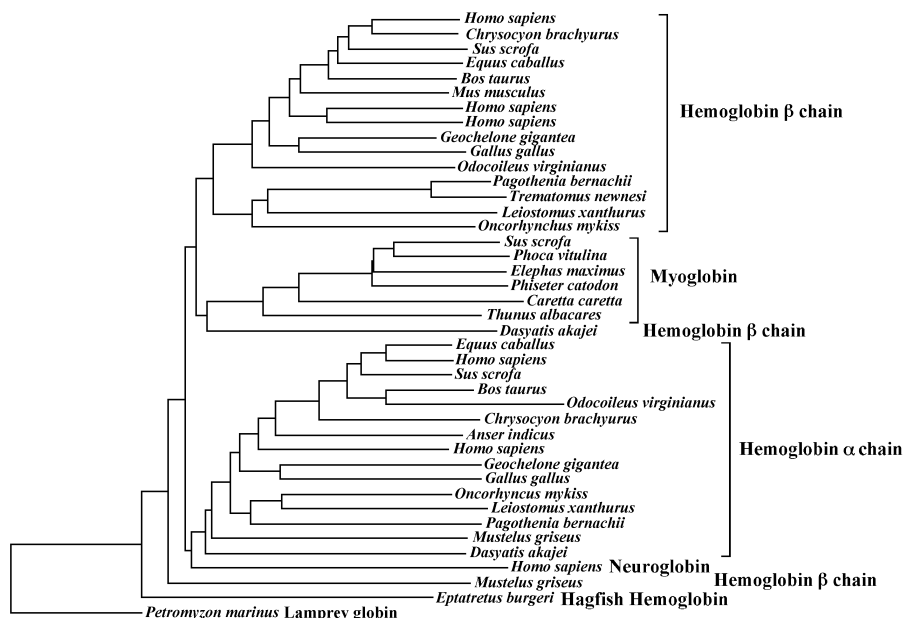
FIGURE 2. Comparison of the SPST and PAM distance matrices.

#### 4. RESULTS

We implemented an algorithm coded in C, called Phyl-SPST, to calculate distances between context trees, as defined by (3.1). The source code and compiled versions for Mac OS X, Linux/Unix and Windows can be downloaded from the site <http://www.ime.usp.br/numec/software/phyl-spst/>.

We applied the Phyl-SPST package to study the similarity between the protein sequences of the globin family of vertebrates. The 41 sequences used in this analysis were obtained from the SCOP database (Andreeva et al.; 2004) and can be found in the supplementary material. The program estimated, for each sequence in this set, a sparse context tree. Then it computed the distance matrix using the  $\beta$ -distance defined by (3.1). In what follows we call this distance the SPST distance. In order to compare our method with an alignment-based distance we used the structure based alignment of the 41 globin sequences of vertebrates present in the PALI database (Gowri et al.; 2003) (alignment available in supplementary material). Then, we applied the algorithm PROTDIST of the Phylip3.65 package (Felsenstein; 2004), with the Dayhoff PAM matrix option, to compute the distance matrix.

(a)



(b)

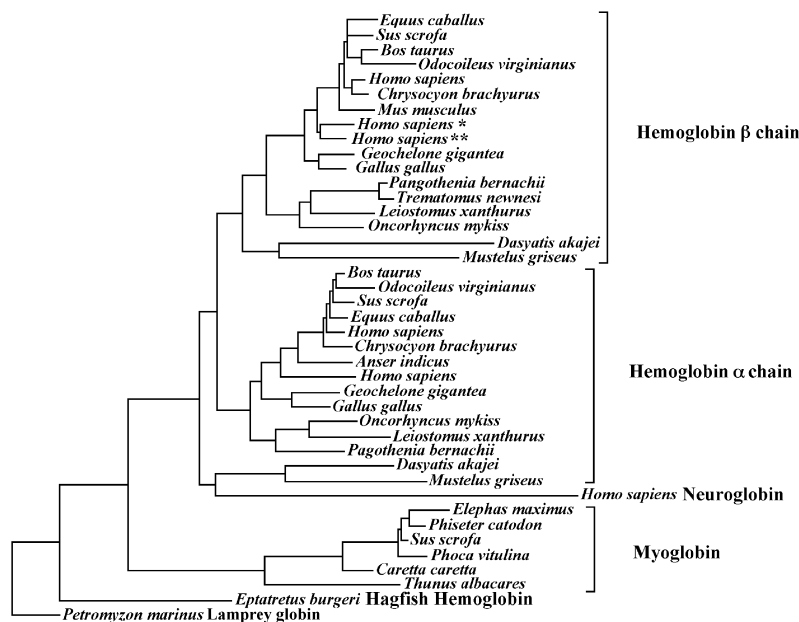


FIGURE 3. Phylogenetic trees made with Neighbor Joining clustering algorithm on SPST distances (a) and on PAM distances (b)

When the PAM and SPST distances are plotted against each other (Fig. 2) a non linear relation is clearly observed. With each distance matrix we reconstructed a phylogenetic tree using the NEIGHBOR and DRAWGRAM algorithms of the Phylip3.65 package. These phylogenetic trees can be seen in Figure 3. In both trees the lamprey globin was used as outgroup.

## 5. DISCUSSION

The dataset we used to verify the potential use of the SPST distances on phylogenetic reconstruction is a vertebrate subset of the globin gene family. This family is one of the first protein families that was characterized (Dayhoff; 1972) and is, perhaps, the most known to date (Vinogradov et al.; 2006). Besides, the vertebrate phylogeny is also well studied and is ground in relatively abundant paleontological, morphological, molecular, and physiological analyses (Cotton and Page; 2002).

The phylogenetic tree shown in Fig. 3(a) proves that in fact the context trees inferred from symbolic sequences (in this case, protein sequences) can offer important evolutionary information of the sequences. This constitutes an original and very promise aspect of the modeling of sequences by variable memory stochastic processes, and it needs to be studied in more details.

The phylogenetic analysis here performed also reflects the overall behavior of the SPST distance. The tree produced with the SPST present larger branches in the most inclusive sequences, and shorter branches in the most basal sequences. With respect to the tree topology, the main differences between them is the placement of the myoglobin cluster, that is closer to the beta chain of hemoglobin in the SPST tree and, in the PAM tree, it is outside of the hemoglobin chain. Other remarkable difference is the placement of the red tail deer (*Odocoileus virginianus*) outside the cluster that contains the mammals, a reptile (*Geochelone gigantea*), and a bird (*Gallus gallus*) in the beta chain cluster of the SPST tree. Although there are minor misplacements in the tree based on PAM distances with respect to the vertebrate and globin traditional phylogenies, it is superior in reconstructing the phylogeny than with the use of SPST distances.

The relationship between the SPST distance and the classical PAM distance of the globin family of vertebrates shows a plateau behavior. The short PAM distances yields larger SPST distances, and the opposite occurs when distances are longer. This may be caused by the bounded nature of the context trees and by the specific form of the distance we propose. Therefore, this analysis shows that small differences in sequences causes enough changes in the context trees to increase the SPST distance between them. It remains yet as an open problem the characterization of the changes produced in the context trees by stationary modifications of the sequences

as mutations, insertions or deletions. We think that these characterizations could help to improve the results shown here. On the other hand, it is also important to define and test other distances over the set of trees to study their specific behaviors and compare them to the one proposed here.

#### ACKNOWLEDGMENTS

This work is part of PRONEX/FAPESP's project *Stochastic behavior, critical phenomena and rhythmic pattern identification in natural languages* (grant number 03/09930-9) and CNPq's project *Stochastic modeling of speech* (grant number 475177/2004-5). During the preparation of this paper F.G.L was supported by a CAPES grant and by a FAPESP fellowship (process 06/56980-0). The authors A.G., H.A.A., and S.R.M. would like to thank the fellowship grants received from CNPq. The authors would like to thank Dr. Eleonora Trajano for comments on the vertebrates phylogenies.

#### REFERENCES

- Andreeva, A., Howorth, D., Brenner, S. E., Hubbard, T. J. P., Chothia, C. and Murzin, A. G. (2004). SCOP database in 2004: refinements integrate structure and sequence family data, *Nucl. Acids Res.* **32**(suppl.1): D226–229.
- Cotton, J. A. and Page, R. (2002). Going nuclear: gene family evolution and vertebrate phylogeny reconciled, *Proc. R. Soc. Lond. B* **269**: 1555–1561.
- Dayhoff, M. (1972). Atlas of protein sequence and structure. National Biomedical Research Foundation, Washington.
- Felsenstein, J. (2004). Phylip (phylogeny inference package) version 3.6. Distributed by the author. Department of Genome Sciences, University of Washington, Seattle.
- Gowri, V. S., Pandit, S. B., Karthik, P. S., Srinivasan, N. and Balaji, S. (2003). Integration of related sequences with protein three-dimensional structural families in an updated version of PALI database, *Nucleic Acids Res.* **31**: 486–488.
- Leonardi, F. (2006). A generalization of the PST algorithm: modeling the sparse nature of protein sequences, *Bioinformatics* **22**(11): 1302–1307.
- Leonardi, F. and Galves, A. (2005). Sequence motif identification and protein family classification using probabilistic trees, *Advances in Bioinformatics and Computational Biology. Proc. BSB 2005.*, Vol. LNBI 3594, pp. 190–193.
- Pham, T. and Zuegg, J. (2004). A probabilistic measure for alignment-free sequence comparison, *Bioinformatics* **20**(18): 3455–3461.
- Rissanen, J. (1983). A universal data compression system, *IEEE Trans. Inform. Theory* **29**(5): 656–664.



- Ron, D., Singer, Y. and Tishby, N. (1996). The power of amnesia: Learning probabilistic automata with variable memory length, *Machine Learning* **25**(2-3): 117–149.
- Simovici, D. and Szymon, J. (2006). A new metric splitting criterion for decision trees, *Journal of Parallel, Emerging and Distributed Computing* **21**(4): 239–256.
- Vinogradov, S. N., Hoogewijs, D., Bailly, X., Arredondo-Peter, R., Gough, J., Dewilde, S., Moens, L. and Vanfleteren, J. R. (2006). A phylogenomic profile of globins, *BMC Evolutionary Biology* **6**: 31–47.
- Wu, T. J., Hsieh, Y. C. and Li, L. A. (2001). Statistical measures of dna dissimilarity under markov chain models of base composition, *Biometrics* **57**: 441–448.

INSTITUTO DE MATEMÁTICA E ESTATÍSTICA, UNIVERSIDADE DE SÃO PAULO., RUA DO MATÃO 1010 CEP 05508-090, SÃO PAULO, SP, BRAZIL.

*E-mail address:* leonardi@ime.usp.br

INSTITUTO DE BIOCÊNCIAS, UNIVERSIDADE DE SÃO PAULO., RUA DO MATÃO, TRAV. 14, N 321 CEP 05508-900, SÃO PAULO, SP, BRAZIL.,

*E-mail address:* srmatiol@ib.usp.br

INSTITUTO DE QUÍMICA, UNIVERSIDADE DE SÃO PAULO., AV. PROF. LINEU PRESTES, 748 CEP 05508-900, SÃO PAULO, SP, BRAZIL.

*E-mail address:* haarmeli@iq.usp.br

INSTITUTO DE MATEMÁTICA E ESTATÍSTICA, UNIVERSIDADE DE SÃO PAULO., RUA DO MATÃO 1010 CEP 05508-090, SÃO PAULO, SP, BRAZIL.,

*E-mail address:* galves@ime.usp.br